



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### An improved pig reference genome sequence to enable pig genetics and genomics research

**Citation for published version:**

Warr, A, Affara, N, Aken, B, Beiki, H, Bickhart, DM, Billis, K, Chow, W, Eory, L, Finlayson, HA, Flicek, P, Girón, CG, Griffin, DK, Hall, R, Hannum, G, Hourlier, T, Howe, K, Hume, DA, Izuogu, O, Kim, K, Koren, S, Liu, H, Manchanda, N, Martin, FJ, Nonneman, DJ, O'Connor, RE, Phillippy, AM, Rohrer, GA, Rosen, BD, Rund, LA, Sargent, CA, Schook, LB, Schroeder, SG, Schwartz, AS, Skinner, BM, Talbot, R, Tseng, E, Tuggle, CK, Watson, M, Smith, TPL & Archibald, AL 2020, 'An improved pig reference genome sequence to enable pig genetics and genomics research', *GigaScience*, vol. 9, no. 6, pp. 1-14.  
<https://doi.org/10.1093/gigascience/giaa051>

**Digital Object Identifier (DOI):**

[10.1093/gigascience/giaa051](https://doi.org/10.1093/gigascience/giaa051)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Publisher's PDF, also known as Version of record

**Published In:**

GigaScience

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



RESEARCH

# An improved pig reference genome sequence to enable pig genetics and genomics research

Amanda Warr <sup>1</sup>, Nabeel Affara<sup>2</sup>, Bronwen Aken <sup>3</sup>, Hamid Beiki <sup>4</sup>, Derek M. Bickhart <sup>5</sup>, Konstantinos Billis <sup>3</sup>, William Chow <sup>6</sup>, Lel Eory <sup>1</sup>, Heather A. Finlayson<sup>1</sup>, Paul Flicek <sup>3</sup>, Carlos G. Girón <sup>3</sup>, Darren K. Griffin <sup>7</sup>, Richard Hall <sup>8</sup>, Greg Hannum<sup>9</sup>, Thibaut Hourlier <sup>3</sup>, Kerstin Howe <sup>6</sup>, David A. Hume <sup>1,10</sup>, Osagie Izuogu <sup>3</sup>, Kristi Kim<sup>8</sup>, Sergey Koren <sup>11</sup>, Haibou Liu<sup>4</sup>, Nancy Manchanda<sup>12</sup>, Fergal J. Martin <sup>3</sup>, Dan J. Nonneman <sup>13</sup>, Rebecca E. O'Connor <sup>7</sup>, Adam M. Phillippy <sup>11</sup>, Gary A. Rohrer <sup>13</sup>, Benjamin D. Rosen <sup>14</sup>, Laurie A. Rund <sup>15</sup>, Carole A. Sargent<sup>2</sup>, Lawrence B. Schook <sup>15</sup>, Steven G. Schroeder <sup>14</sup>, Ariel S. Schwartz<sup>9</sup>, Ben M. Skinner <sup>2</sup>, Richard Talbot<sup>16</sup>, Elizabeth Tseng <sup>8</sup>, Christopher K. Tuggle <sup>4,12</sup>, Mick Watson <sup>1</sup>, Timothy P. L. Smith <sup>13,\*</sup> and Alan L. Archibald <sup>1,\*</sup>

<sup>1</sup>The Roslin Institute and Royal (Dick) School of Veterinary Studies, The University of Edinburgh, Easter Bush Campus, Midlothian EH25 9RG, UK; <sup>2</sup>Department of Pathology, University of Cambridge, Tennis Court Road, Cambridge CB2 1QP, UK; <sup>3</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton CB10 1SD, UK; <sup>4</sup>Department of Animal Science, 2255 Kildee Hall, Iowa State University, Ames, IA 50011-3150, USA; <sup>5</sup>Dairy Forage Research Center, USDA-ARS, 1925 Linden Drive, Madison, WI 53706, USA; <sup>6</sup>Wellcome Sanger Institute, Wellcome Genome Campus, Cambridge CB10 1SA, UK; <sup>7</sup>School of Biosciences, University of Kent, Giles Lane, Canterbury CT2 7NJ, UK; <sup>8</sup>Pacific Biosciences, 1305 O'Brien Drive, Menlo Park, CA 94025, USA; <sup>9</sup>Denovium Inc., San Diego, CA, USA; <sup>10</sup>Mater Research Institute-University of Queensland, Translational Research Institute, Brisbane QLD 4104, Australia; <sup>11</sup>Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, National Institutes of Health, 9000 Rockville Pike, Bethesda, MD 20892, USA; <sup>12</sup>Bioinformatics and Computational Biology Program, Iowa State University, 2014 Molecular Biology Building, Ames, IA 50011, USA; <sup>13</sup>USDA-ARS U.S. Meat Animal Research Center, 844 Road 313, Clay Center, NE 68933, USA; <sup>14</sup>Animal Genomics and Improvement Laboratory, USDA-ARS, 10300 Baltimore Avenue, Beltsville, MD 20705-2350, USA;

Received: 28 October 2019; Revised: 12 March 2020; Accepted: 22 April 2020

© The Author(s) 2020. Published by Oxford University Press. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

<sup>15</sup>Department of Animal Sciences, University of Illinois, 1201 West Gregory Drive, Urbana, IL 61801, USA and

<sup>16</sup>Edinburgh Genomics, University of Edinburgh, Charlotte Auerbach Road, Edinburgh EH9 3FL, UK

\*Correspondence address. Alan L. Archibald, The Roslin Institute and Royal (Dick) School of Veterinary Studies, The University of Edinburgh, Easter Bush Campus, Midlothian EH25 9RG, UK. E-mail: [alan.archibald@roslin.ed.ac.uk](mailto:alan.archibald@roslin.ed.ac.uk)  <http://orcid.org/0000-0001-9213-1830>; Timothy P.L. Smith, U.S. Meat Animal Research Center, USDA-ARS, 844 Road 313, Clay Center, NE 68933, USA. E-mail: [tim.smith2@usda.gov](mailto:tim.smith2@usda.gov)  <http://orcid.org/0000-0003-1611-6828>

## Abstract

**Background:** The domestic pig (*Sus scrofa*) is important both as a food source and as a biomedical model given its similarity in size, anatomy, physiology, metabolism, pathology, and pharmacology to humans. The draft reference genome (Sscrofa10.2) of a purebred Duroc female pig established using older clone-based sequencing methods was incomplete, and unresolved redundancies, short-range order and orientation errors, and associated misassembled genes limited its utility.

**Results:** We present 2 annotated highly contiguous chromosome-level genome assemblies created with more recent long-read technologies and a whole-genome shotgun strategy, 1 for the same Duroc female (Sscrofa11.1) and 1 for an outbred, composite-breed male (USMARCv1.0). Both assemblies are of substantially higher (>90-fold) continuity and accuracy than Sscrofa10.2. **Conclusions:** These highly contiguous assemblies plus annotation of a further 11 short-read assemblies provide an unprecedented view of the genetic make-up of this important agricultural and biomedical model species. We propose that the improved Duroc assembly (Sscrofa11.1) become the reference genome for genomic research in pigs.

**Keywords:** pig genomes; reference assembly; pig; genome annotation

## Background

High-quality, richly annotated reference genome sequences are key resources and provide important frameworks for the discovery and analysis of genetic variation and for linking genotypes to function. In farmed animal species such as the domestic pig (*Sus scrofa*, NCBI:txid9823) genome sequences have been integral to the discovery of molecular genetic variants and the development of single-nucleotide polymorphism (SNP) chips [1] and enabled efforts to dissect the genetic control of complex traits, such as growth, feed conversion, body composition, reproduction, behaviour, and responses to infectious diseases [2].

Genome sequences are an essential resource not only for enabling research but also for applications in the life sciences. Genomic selection, in which associations between thousands of SNPs and trait variation as established in a phenotyped training population are used to choose amongst selection candidates for which there are SNP data but no phenotypes, has delivered genomics-enabled genetic improvement in farmed animals [3] and plants. From its initial successful application in dairy cattle breeding, genomic selection is now being used in many sectors within animal and plant breeding, including by leading pig breeding companies [4, 5].

The domestic pig (*S. scrofa*) has importance not only as a source of animal protein but also as a biomedical model. The choice of the optimal animal model species for pharmacological or toxicology studies can be informed by knowledge of the genome and gene content of the candidate species including pigs [6]. A high quality, richly annotated genome sequence is also essential when using gene editing technologies to engineer improved animal models for research or as sources of cells and tissue for xenotransplantation and potentially for improved productivity [7, 8].

The highly continuous pig genome sequences reported here are built upon a quarter of a century of effort by the global pig genetics and genomics research community including the development of recombination and radiation hybrid (RH) maps [9, 10], cytogenetic and bacterial artificial chromosome (BAC) physical maps [11, 12], and a draft reference genome sequence [13].

The previously published draft pig reference genome sequence (Sscrofa10.2), developed under the auspices of the Swine Genome Sequencing Consortium (SGSC), has a number of consequential deficiencies [14–17]. The BAC-by-BAC hierarchical shotgun sequence approach [18] using Sanger sequencing technology can yield a high quality genome sequence as demonstrated by the public Human Genome Project. However, with a fraction of the financial resources of the Human Genome Project, the resulting draft pig genome sequence comprised an assembly, in which long-range order and orientation is good, but the order and orientation of sequence contigs within many BAC clones was poorly supported and the sequence redundancy between overlapping sequenced BAC clones was often not resolved. Moreover, ~10% of the pig genome, including some important genes, was not represented (e.g., CD163) or incompletely represented (e.g., IGF2) in the assembly [19]. Whilst the BAC clones represent an invaluable resource for targeted sequence improvement and gap closure as demonstrated for chromosome X (SSCX) [20], a clone-by-clone approach to sequence improvement is expensive notwithstanding the reduced cost of sequencing with next-generation technologies.

The dramatically reduced cost of whole-genome shotgun sequencing using Illumina short-read technology has facilitated the sequencing of several hundred pig genomes [17, 21, 22]. Whilst a few of these additional pig genomes have been assembled to contig level, most of these genome sequences have simply been aligned to the reference and used as a resource for variant discovery.

The increased capability and reduced cost of third-generation long-read sequencing technology as delivered by Pacific Biosciences (PacBio) and Oxford Nanopore platforms have created the opportunity to generate the data from which to build highly contiguous genome sequences as illustrated recently for cattle [23, 24]. Here we describe the use of PacBio long-read technology to establish highly continuous pig genome sequences that provide substantially improved resources for pig genetics and genomics research and applications.

**Table 1:** Assembly statistics

Statistic	Sscrofa10.2	Sscrofa11	Sscrofa11.1	USMARCv1.0	GRCh38.p13
Total sequence length	2,808,525,991	2,456,768,445	2,501,912,388	2,755,438,182	3,099,706,404
Total ungapped length	2,519,152,092	2,454,899,091	2,472,047,747	2,623,130,238	2,948,583,725
No. of scaffolds	9,906	626	706	14,157	472
Gaps between scaffolds	5,323	24	93	0	349
No. of unplaced scaffolds	4,562	583	583	14,136	126
Scaffold N50	576,008	88,231,837	88,231,837	131,458,098	67,794,873
Scaffold L50	1,303	9	9	9	16
No. of unspanned gaps	5,323	24	93	0	349
No. of spanned gaps	233,116	79	413	661	526
No. of contigs	243,021	705	1,118	14,818	998
Contig N50	69,503	48,231,277	48,231,277	6,372,407	57,879,411
Contig L50	8,632	15	15	104	18
No. of chromosomes*	*21	19	*21	*21	24

Summary statistics for assembled pig genome sequences and comparison with current human reference genome (source: NCBI, <https://www.ncbi.nlm.nih.gov/assembly/>). \*Includes mitochondrial genome.

## Results

Two individual pigs were sequenced independently: (i) TJ Tabasco (Duroc 2–14), i.e., the sow that was the primary source of DNA for the published draft genome sequence (Sscrofa10.2) [13] and (ii) MARC1423004, which was a crossbred barrow (i.e., castrated male pig) from a composite population (approximately one-half Landrace, one-quarter Duroc, and one-quarter Yorkshire) at the United States Department of Agriculture (USDA) Meat Animal Research Center. The former allowed us to build upon the earlier draft genome sequence, exploit the associated CHORI-242 BAC library resource [25], and evaluate the improvements achieved by comparison with Sscrofa10.2. The latter allowed us to assess the relative efficacy of a simpler whole-genome shotgun sequencing and Chicago Hi-Rise scaffolding strategy [26]. This second assembly also provided data for the Y chromosome and supported comparison of haplotypes between individuals. In addition, full-length transcript sequences were collected for multiple tissues from the MARC1423004 animal and used in annotating both genomes.

### Sscrofa11.1 assembly

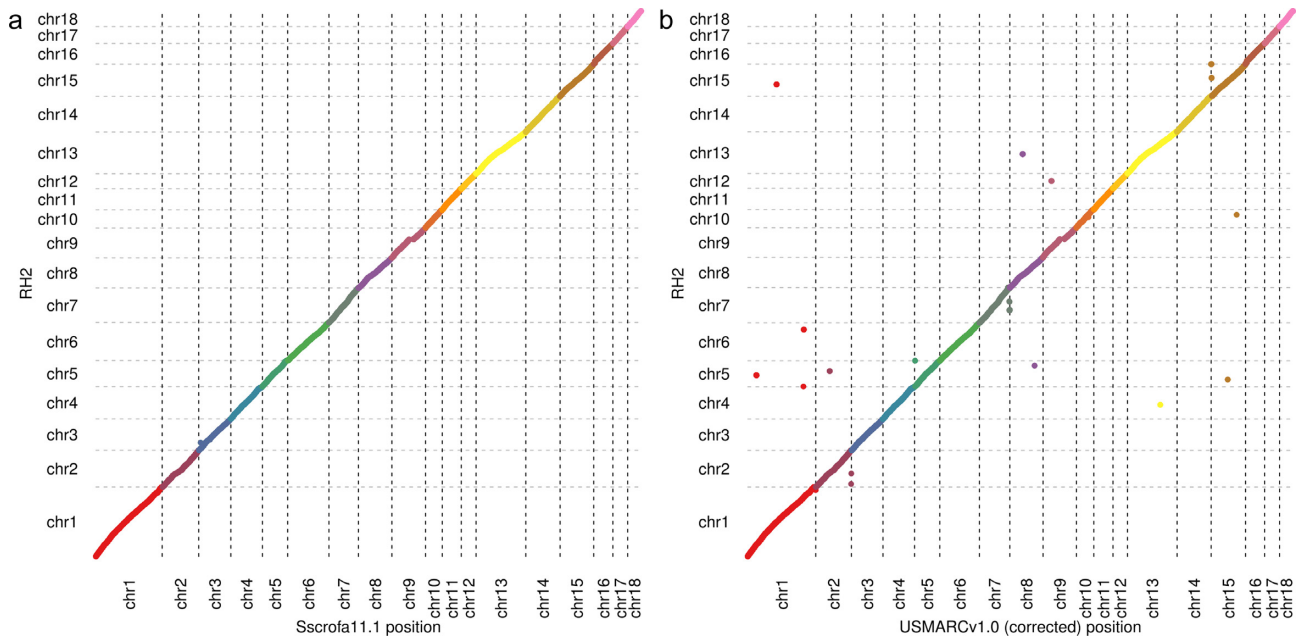
Approximately 65-fold coverage (176 Gb) of the genome of TJ Tabasco (Duroc 2–14) was generated using PacBio single-molecule real-time (SMRT) sequencing technology. A total of 213 SMRT cells produced 12,328,735 subreads of mean length 14,270 bp and with a read N50 of 19,786 bp (Table S1). Reads were corrected and assembled using Falcon (v.0.4.0) [27], achieving a minimum corrected read cut-off of 13 kb that provided 19-fold genome coverage for input, resulting in an initial assembly comprising 3,206 contigs with a contig N50 of 14.5 Mb.

The contigs were mapped to the previous draft assembly (Sscrofa10.2) using Nucmer [28]. The long-range order of the Sscrofa10.2 assembly was based on fingerprint contig [12] and RH physical maps with assignments to chromosomes based on fluorescence in situ hybridization (FISH) data. This alignment of Sscrofa10.2 and the contigs from the initial Falcon assembly of the PacBio data provided draft scaffolds that were tested for consistency with paired BAC and fosmid end sequences and the RH map [9]. The draft scaffolds also provided a framework for gap closure using PBjelly [29], or finished quality Sanger sequence data generated from CHORI-242 BAC clones from earlier work [13, 20].

Remaining gaps between contigs within scaffolds, and between scaffolds predicted to be adjacent on the basis of other available data, were targeted for gap filling with a combination of unplaced contigs and previously sequenced BACs, or by identification and sequencing of BAC clones predicted from their end sequences to span the gaps. The combination of methods filled 2,501 gaps and reduced the number of contigs in the assembly from 3,206 to 705. The assembly, Sscrofa11 (GCA.000003025.5), had a final contig N50 of 48.2 Mb, only 103 gaps in the sequences assigned to chromosomes, and only 583 remaining unplaced contigs (Table 1). Two acrocentric chromosomes (SSC16, SSC18) were each represented by single, unbroken contigs. The SSC18 assembly also includes centromeric and telomeric repeats (Tables S2 and S3; Figs S1 and S2), albeit the former probably represent a collapsed version of the true centromere. The reference genome assembly was completed by adding Y chromosome sequences from other sources (GCA.900119615.2) [20] because TJ Tabasco (Duroc 2–14) was female. The resulting reference genome sequence was termed Sscrofa11.1 and deposited in the public sequence databases (GCA.000003025.6) (Table 1).

The medium- to long-range order and orientation of the Sscrofa11.1 assembly was assessed by comparison with an existing RH map [9]. The comparison strongly supported the overall accuracy of the assembly (Fig. 1a), despite the fact that the RH map was prepared from a cell line of a different individual. There is 1 major disagreement between the RH map and the assembly on chromosome 3, which will need further investigation. The only other substantial disagreement on chromosome 9 is explained by a gap in the RH map [9]. The assignment and orientation of the Sscrofa11.1 scaffolds to chromosomes was confirmed with FISH of BAC clones (Table S4, Fig. S3). The Sscrofa11.1 and USMARCv1.0 assemblies were searched using BLAST [30] with sequences derived from the BAC clones that had been used as probes for the FISH analyses. For most BAC clones these sequences were BAC end sequences [12], but in some cases these sequences were incomplete or complete BAC clone sequences [13, 20]. The links between the genome sequence and the BAC clones used in cytogenetic analyses by FISH are summarized in Table S4. The FISH results indicate areas where future assemblies might be improved. For example, the Sscrofa11.1 unplaced scaffolds contig1206 and contig1914 may contain sequences that could be added to the ends of the long arms of SSC1 and SSC7, respectively.





**Figure 1:** Assemblies and radiation hybrid (RH) map alignments. Plots illustrating co-linearity between RH map and (a) Sscrofa11.1 and (b) USMARCv1.0 assemblies (autosomes only).

The quality of the Sscrofa11 assembly, which corresponds to Sscrofa11.1 after the exclusion of SSCY, was assessed as described previously for the existing Sanger sequence-based draft assembly (Sscrofa10.2) [14]. Alignments of Illumina sequence reads from the same female pig were used to identify regions of low quality (LQ; regions with high GC normalized coverage, prevalence of improperly paired reads, and prevalence of reads with improper insert sizes) or low coverage (LC; regions with low GC normalized coverage) (Table 2). The analysis confirms that Sscrofa11 represents a substantial improvement over the Sscrofa10.2 draft assembly. For example, the low-quality low-coverage (LQLC) proportion of the genome sequence has decreased from 33.1% to 16.3% when repetitive sequence is not masked and to 1.6% when repeats are masked prior to read alignment. The remaining LQLC segments of Sscrofa11 have a mean GC content of 61.6%. Thus, these regions may represent sequence where short-read coverage is low as a result of the known systematic bias of the short-read platform against extreme GC content sequences, rather than deficiencies of the assembly.

The Sscrofa11.1 assembly was also assessed visually using gEVAL [31]. The improvement in short-range order and orientation as revealed by alignments with isogenic BAC and fosmid end sequences is illustrated for a particularly poor region of Sscrofa10.2 on chromosome 12 (Fig. S4). The problems in this area of Sscrofa10.2 arose from failures to order and orient the sequence contigs and resolve the redundancies between these sequence contigs within BAC clone CH242-147O24 (ENA: FP102566.2). The improved contiguity in Sscrofa11.1 not only resolves these local order and orientation errors but also facilitates the annotation of a complete gene model for the ABR locus. Further examples of comparisons of Sscrofa10.2 and Sscrofa11.1 reveal improvements in contiguity, local order and orientation, and gene models (Figs S5–S7).

### USMARCv1.0 assembly

Approximately 65-fold coverage of the genome of the MARC1423004 barrow was generated on a PacBio RSII instrument. The sequence was collected during the transition from P5/C3 to P6/C4 chemistry, with approximately equal numbers of subreads from each chemistry. A total of 199 cells of P5/C3 chemistry produced 95.3 Gb of sequence with mean subread length of 5.1 kb and subread N50 of 8.2 kb. A total of 127 cells of P6/C4 chemistry produced 91.6 Gb of sequence with mean subread length 6.5 kb and subread N50 of 10.3 kb, resulting in an overall mean subread length, including data from both chemistries, of 6.4 kb. The reads were assembled using Celera Assembler 8.3rc2 [32] and Falcon [27]. The resulting assemblies were compared, and the Celera Assembler result was selected on the basis of better agreement with a Dovetail Chicago® library [26] (i.e., there was a lower proportion of conflicting links between read pairs from the Chicago library) and was used to create a scaffolded assembly with the HiRise™ scaffolder consisting of 14,818 contigs with a contig N50 of 6.372 Mb (GenBank accession GCA.002844635.1; Table 1). The USMARCv1.0 scaffolds were therefore completely independent of the existing Sscrofa10.2 or new Sscrofa11.1 assemblies, and they can act as supporting evidence where they agree with those assemblies. However, chromosome assignment of the scaffolds was performed by alignment to Sscrofa10.2 and does not constitute independent confirmation of this ordering. The assignment of these scaffolds to individual chromosomes was confirmed post hoc by FISH analysis as described for Sscrofa11.1 above. The FISH analysis revealed that several of these chromosome assemblies (SSC1, 5, 6–11, 13–16) are inverted with respect to the cytogenetic convention for pig chromosome (Table S4; Figs S3 and S8–S10). After correcting the orientation of these inverted scaffolds, there is good agreement between the USMARCv1.0 assembly and the RH map [9] (Fig. 1b, Table S5).

**Table 2:** Summary of quality statistics for SSC1–18, SSCX

Statistic	Bases, Sscrofa11	Sscrofa11	% Genome	Sscrofa10.2
High coverage	119,341,205	4.9		2.6
LC	185,385,536	7.5		26.6
Low proportion properly paired	95,508,007	3.9		5.0
High proportion large inserts	40,835,320	1.7		1.5
High proportion small inserts	114,793,298	4.7		4.0
LQ	284,838,040	11.6		13.9
Total LQLC	399,927,747	16.3		33.1
LQLC windows that do not intersect RepeatMasker regions	39,918,551	1.6		

Quality measures and terms as defined [14]. LC: low coverage; LQ: low quality.

### Sscrofa11.1 and USMARCv1.0 are co-linear

The alignment of the 2 PacBio assemblies reveals a high degree of agreement and co-linearity, after correction of the inversions of several USMARCv1.0 chromosome assemblies (Fig. S11). The agreement between the Sscrofa11.1 and USMARCv1.0 assemblies is also evident in comparisons of specific loci (Figs S5–S7) although with some differences (e.g., Fig. S6). The whole-genome alignment of Sscrofa11.1 and USMARCv1.0 (Fig. S11) masks some inconsistencies that are evident when the alignments are viewed on a single chromosome-by-chromosome basis (Figs S8–S10). It remains to be determined whether the small differences between the assemblies represent errors in the assemblies or true structural variation between the 2 individuals (see discussion of the *ERLIN1* locus below).

Pairwise comparisons amongst the Sscrofa10.2, Sscrofa11.1, and USMARCv1.0 assemblies using the Assemblytics tools [33] revealed a peak of insertions and deletion with sizes of ~300 bp (Figs S12a–S12c). We assume that these correspond to short interspersed nuclear elements. Both the Sscrofa11.1 and USMARCv1.0 assemblies have more differences against Sscrofa10.2 (33,347 and 44,023, respectively) than against each other (28,733). This is despite the fact that Sscrofa11.1 and Sscrofa10.2 represent the same pig genome. While some differences between Sscrofa10.2 and Sscrofa11.1 may be due to differences in which haplotype has been captured in the assembly, the reduction in LQ and LC regions and the dramatic decrease in differences versus USMARCv1.0 lead us to conclude that the majority are improvements in the Sscrofa11.1 assembly. The differences between Sscrofa11.1 and USMARCv1.0 will represent a mix of true structural differences and assembly errors that will require further research to resolve. The Sscrofa11.1 and USMARCv1.0 assemblies were also compared with 11 Illumina short-read assemblies [17] (Table S6).

### Repetitive sequences, centromeres, and telomeres

The repetitive sequence content of Sscrofa11.1 and USMARCv1.0 was identified and characterized. These analyses allowed the identification of centromeres and telomeres for several chromosomes. The previous reference genome (Sscrofa10.2) that was established from Sanger sequence data and a minipig genome (minipig.v1.0, GCA.000325925.2) that was established from Illumina short-read sequence data were also included for comparison. The numbers of the different repeat classes and the average mapped lengths of the repetitive elements identified in these 4 pig genome assemblies are summarized in Figs S13 and S14, respectively.

Putative telomeres were identified at the proximal ends of Sscrofa11.1 chromosome assemblies of SSC2, SSC3, SSC6, SSC8, SSC9, SSC14, SSC15, SSC18, and SSCX (Fig. S1; Table S2). Putative centromeres were identified in the expected locations in the Sscrofa11.1 chromosome assemblies for SSC1–7, SSC9, SSC13, and SSC18 (Fig. S2, Table S3). For the chromosome assemblies of each of SSC8, SSC11, and SSC15, 2 regions harbouring centromeric repeats were identified. Pig chromosomes SSC1–12 plus SSCX and SSCY are all metacentric, whilst chromosomes SSC13–18 are acrocentric. The putative centromeric repeats on SSC17 do not map to the expected end of the chromosome assembly.

### Completeness of the assemblies

The Sscrofa11.1 and USMARCv1.0 assemblies were assessed for completeness using 2 tools, BUSCO [34] and Cogent [35]. BUSCO uses a database of expected gene content based on near-universal single-copy orthologs from species with genomic data, while Cogent uses transcriptome data from the organism being sequenced and therefore provides an organism-specific view of genome completeness. BUSCO analysis suggests that both new assemblies are highly complete, with 93.8% and 93.1% of BUSCOs complete for Sscrofa11.1 and USMARCv1.0, respectively, a marked improvement on the 80.9% complete in Sscrofa10.2 and comparable to the human and mouse reference genome assemblies (Table S7).

Cogent is a tool that identifies gene families and reconstructs the coding genome using full-length, high-quality (HQ) transcriptome data without a reference genome and can be used to check assemblies for the presence of these known coding sequences [35]. PacBio transcriptome (Iso-Seq) data consisting of HQ isoform sequences from 7 tissues (diaphragm, hypothalamus, liver, skeletal muscle [longissimus dorsi], small intestine, spleen, and thymus) [36] from the pig whose DNA was used as the source for the USMARCv1.0 assembly were pooled together for Cogent analysis. Cogent partitioned 276,196 HQ isoform sequences into 30,628 gene families, of which 61% had  $\geq 2$  distinct transcript isoforms. Cogent then performed reconstruction on the 18,708 partitions. For each partition, Cogent attempts to reconstruct coding “contigs” that represent the ordered concatenation of transcribed exons as supported by the isoform sequences. The reconstructed contigs were then mapped back to Sscrofa11.1, and contigs that could not be mapped or map to  $>1$  position were individually examined. There were 5 genes that were present in the Iso-Seq data but missing in the Sscrofa11.1 assembly. In each of these 5 cases, a Cogent partition (which consists of  $\geq 2$  transcript isoforms of the same gene, often from multiple tissues) exists in which the predicted transcript does not

**Table 3:** Annotation statistics for Ensembl annotation of pig (Sscrofa10.2, Sscrofa11.1, USMARCv1.0), human (GRCh38.p13), and mouse (GRCm38.p6) assemblies

Statistic	Sscrofa10.2 (Release 89)	Sscrofa11.1 (Release 98)	USMARCv1.0 (Release 97)	GRCh38.p13 (Release 98)	GRCm38.p6 (Release 98)
Coding genes	21,630 (incl 10 RT)	21,301	21,535	20,444 (incl 667 RT)	22,508 (incl 270 RT)
Non-coding genes	3,124	8,971	6,113	23,949	16,078
Small non-coding genes	2,804	2,156	2,427	4,871	5,531
Long non-coding genes	135 (incl 1 RT)	6,798	3,307	16,857 (incl 304 RT)	9,985 (incl 75 RT)
Miscellaneous non-coding genes	185	17	379	2,221	562
Pseudogenes	568	1,626	674	15,214 (incl 8 RT)	13,597 (incl 4 RT)
Gene transcripts	30,585	63,041	58,692	227,530	142,446
Genscan gene predictions	52,372	46,573	152,168	51,756	57,381
Short variants	60,389,665	64,310,125		665,834,144	83,761,978
Structural variants	224,038	224,038		6,013,113	791,878

Incl: including; RT: read through.

align back to Sscrofa11.1. NCBI-BLASTN of the isoforms from the partitions revealed them to have near-perfect hits with existing annotations for *CHAMP1*, *ERLIN1*, *IL1RN*, *MB*, and *PSD4* for other species.

*ERLIN1* is missing from its predicted location on SSC14 between the *CHUK* and *CPN1* genes in Sscrofa11.1. There is good support for the Sscrofa11.1 assembly in the region from the BAC end sequence alignments, suggesting that this area may represent a true haplotype. Indeed, a copy number variant nsv1302227 has been mapped to this location on SSC14 [37] and the *ERLIN1* gene sequences present in BAC clone CH242-513L2 (ENA: CT868715.3) were incorporated into the earlier Sscrofa10.2 assembly. However, an alternative haplotype containing *ERLIN1* was not found in any of the assembled contigs from Falcon and this will require further investigation. The *ERLIN1* locus is present on SSC14 in the USMARCv1.0 assembly (30,107,816–30,143,074; note that the USMARCv1.0 assembly of SSC14 is inverted relative to Sscrofa11.1). Of 11 short-read pig genome assemblies [17] that have been annotated with the Ensembl pipeline (Ensembl release 98, September 2019) [38, 39], *ERLIN1* sequences are present in the expected genomic context in all 11 genome assemblies. The fact that the *ERLIN1* gene is located at the end of a contig in 8 of these short-read assemblies suggests that this region of the pig genome presents difficulties for sequencing and assembly and the absence of *ERLIN1* in Sscrofa11.1 is more likely to be an assembly error.

The other 4 genes are annotated in neither Sscrofa10.2 nor Sscrofa11.1. Two of these genes, *IL1RN* and *PSD4*, are present in the original Falcon contigs; however, they were trimmed off during the contig QC stage because of apparent abnormal Illumina, BAC, and fosmid mapping in the region, which was likely caused by the repetitive nature of their expected location on chromosome 3 where a gap is present. The *IL1RN* and *PSD4* genes are present in USMARCv1.0, albeit their location is anomalous, and are also present in the 11 short-read assemblies [17]. *CHAMP1* (ENSSSCG00070014091) is present in the USMARCv1.0 assembly in the sub-telomeric region of the q-arm, after correction of the inversion of the USMARCv1.0 scaffold, and is also present in all 11 short-read assemblies [17]. After correction of the orientation of the USMARCv1.0 chromosome 11 scaffold there is a small inversion of the distal 1.07 Mb relative to the Sscrofa11.1 assembly; this region harbours the *CHAMP1* gene. The orientation of the Sscrofa11.1 chromosome 11 assembly in this region is consistent with the predictions of the human-pig comparative map [40].

The myoglobin gene (*MB*) is present in the expected location in the USMARCv1.0 assembly flanked by *RASD2* and *RBFOX2*. Partial *MB* sequences are present distal to *RBFOX2* on chromosome 5 in the Sscrofa11.1 assembly. Because there is no gap here in the Sscrofa11.1 assembly it is likely that the incomplete *MB* is a result of a misassembly in this region. This interpretation is supported by a break in the pairs of BAC and fosmid end sequences that map to this region of the Sscrofa11.1 assembly. Some of the expected gene content missing from this region of the Sscrofa11.1 chromosome 5 assembly, including *RASD2*, *HMOX1*, and *LARGE1*, is present on an unplaced scaffold (AEMK02000361.1). Cogent analysis also identified 2 cases of potential fragmentation in the Sscrofa11.1 genome assembly that resulted in the isoforms being mapped to 2 separate loci, although these will require further investigation. In summary, the BUSCO and Cogent analyses indicate that the Sscrofa11.1 assembly captures a very high proportion of the expressed elements of the genome.

### Improved annotation

Annotation of Sscrofa11.1 was carried out with the Ensembl annotation pipeline and released via the Ensembl Genome Browser (Ensembl release 90, August 2017) [38, 41]. Statistics for the annotation as updated in June 2019 (Ensembl release 98, September 2019) are listed in Table 3. This annotation is more complete than that of Sscrofa10.2 and includes fewer fragmented genes and pseudogenes.

The annotation pipeline used extensive short-read RNA-sequencing (RNA-Seq) data from 27 tissues and long-read PacBio Iso-Seq data from 9 adult tissues. This provided an unprecedented window into the pig transcriptome and allowed for not only an improvement to the main gene set but also the generation of tissue-specific gene tracks from each tissue sample. The use of Iso-Seq data also improved the annotation of untranslated regions because they represent transcripts sequenced across their full length from the polyA tract.

In addition to improved gene models, annotation of the Sscrofa11.1 assembly provides a more complete view of the porcine transcriptome than annotation of the previous assembly (Sscrofa10.2; Ensembl releases 67–89, May 2012 through May 2017) [42], with increases in the numbers of transcripts annotated (Table 3). However, the number of annotated transcripts remains lower than in the human and mouse genomes. The annotation of the human and mouse genomes and in particular the



gene content and encoded transcripts has been more thorough as a result of extensive manual annotation.

Efforts were made to annotate important classes of genes, in particular immunoglobulins and olfactory receptors. For these genes, sequences were downloaded from specialist databases and the literature to capture as much detail as possible (see supplementary information, section 2 annotation, for more details).

These improvements in terms of the resulting annotation were evident in the results of the comparative genomics analyses run on the gene set. The previous annotation had 12,919 one-to-one orthologs with human, while the new annotation of the Sscrofa11.1 assembly has 15,544. Similarly, in terms of conservation of synteny, the previous annotation had 11,661 genes with high-confidence gene order conservation scores, while the new annotation has 15,958. There was also a large reduction in terms of genes that were either abnormally short or split when compared to their orthologs in the new annotation.

The Sscrofa11.1 assembly has also been annotated using the NCBI pipeline [43]. We have compared these 2 annotations. The Ensembl and NCBI annotations of Sscrofa11.1 are broadly similar (Table S8). There are 17,676 protein-coding genes and 1,700 non-coding genes in common. However, 540 of the genes annotated as protein-coding by Ensembl are annotated as non-coding or pseudogenes by NCBI and 227 genes annotated as non-coding by NCBI are annotated as protein-coding (215) or as pseudogenes (12) by Ensembl. The NCBI RefSeq annotation can be visualized in the Ensembl Genome Browser by loading the RefSeq GFF3 track and the annotations compared at the individual locus level. Similarly, the Ensembl annotated genes can be visualized in the NCBI Genome Browser. Despite considerable investment there are also differences in the Ensembl and NCBI annotation of the human reference genome sequence, with 20,444 and 19,755 protein-coding genes on the primary assembly, respectively. The MANE (Matched Annotation from NCBI and EMBL-EBI) project was launched to resolve these differences and identify a matched representative transcript for each human protein-coding gene [44]. To date a MANE transcript has been identified for 12,985 genes.

We have also annotated the USMARCv1.0 assembly using the Ensembl pipeline [38], and this annotation was released via the Ensembl Genome Browser (Ensembl release 97, July 2019) [39] (see Table 3 for summary statistics). More recently, we have annotated a further 11 short-read pig genome assemblies [17] (Ensembl release 98, September 2019) [39]; see Tables S6 and S11 for summary statistics for the assemblies and annotation, respectively.

### SNP chip probes mapped to assemblies

The probes from 4 commercial SNP chips were mapped to the Sscrofa10.2, Sscrofa11.1, and USMARCv1.0 assemblies. We identified 1,709, 56, and 224 markers on the PorcineSNP60, GGP LD, and 80 K commercial chips that were previously unmapped and now have coordinates on the Sscrofa11.1 reference (Table S9). These newly mapped markers can now be imputed into a cross-platform, common set of SNP markers for use in genomic selection. Additionally, we have identified areas of the genome that are poorly tracked by the current set of commercial SNP markers. The previous Sscrofa10.2 reference had a mean (SD) marker spacing of 3.57 (26.5) kb with markers from 4 commercial genotyping arrays. We found this to be an underestimate of the actual distance between markers because the Sscrofa11.1 reference coordinates consisted of a mean (SD) of 3.91 (14.9) kb between the

same set of markers. We also found a region of 2.56 Mb that is currently devoid of suitable markers on the new reference.

A Spearman rank order ( $\rho$ ) value was calculated for each assembly (alternative hypothesis:  $\rho = 0$ ;  $P < 2.2 \times 10^{-16}$ ): Sscrofa10.2: 0.88464; Sscrofa11.1: 0.88890; USMARCv1.0: 0.81260. This rank order comparison was estimated by ordering all of the SNP probes from all chips by their listed manifest coordinates against their relative order in each assembly (with chromosomes ordered by karyotype). Any unmapped markers in an assembly were penalized by giving the marker a “-1” rank in the assembly ranking order.

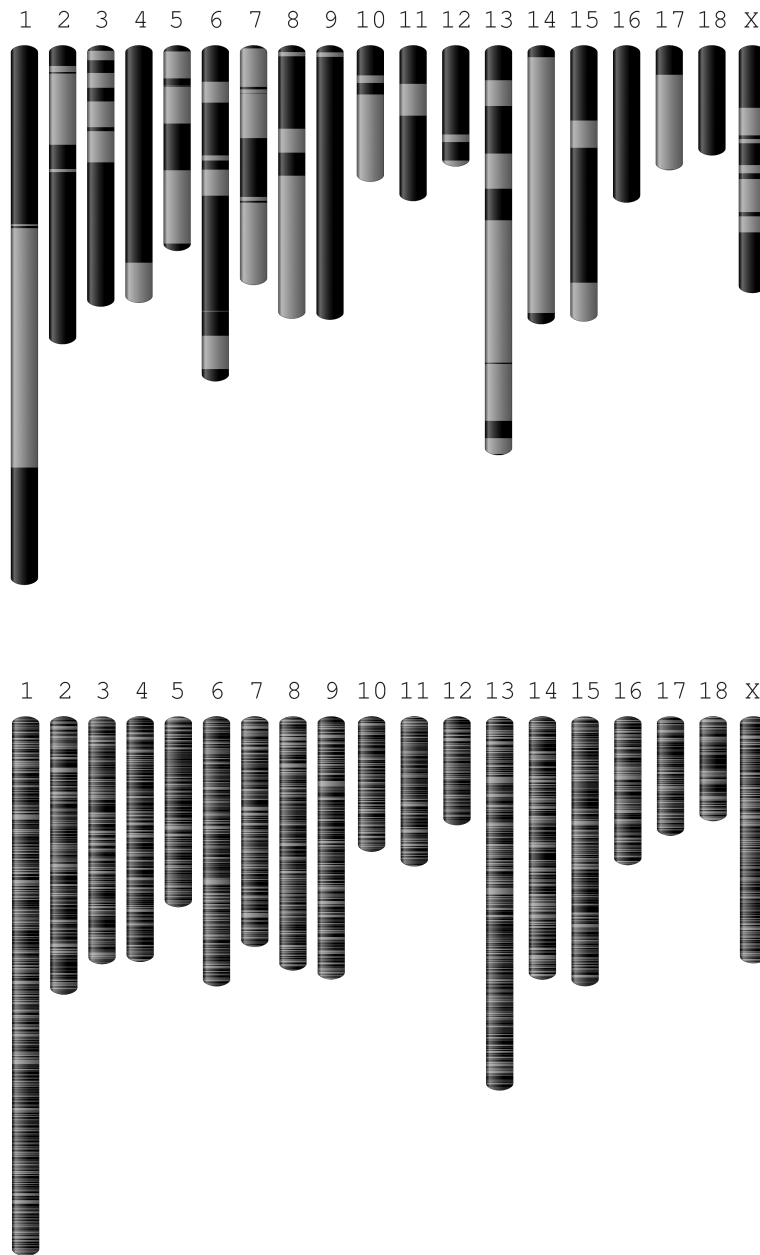
To examine the general linear order of placed markers on each assembly, the marker rank order (y-axis; used above in the Spearman rank order test) was plotted against the probe rank order on the manifest file (x-axis) (Fig. S15). The analyses revealed some interesting artefacts that suggest that the SNP manifest coordinates for the porcine 60 K SNP chip are still derived from an obsolete (Sscrofa9) reference in contrast to all other manifests (Sscrofa10.2). Also, it confirms that several of the USMARCv1.0 chromosome scaffolds are inverted with respect to the canonical orientation of pig chromosomes. The large band of points at the top of the plot corresponds to marker mappings on the unplaced contigs of each assembly. These unplaced contigs often correspond to assemblies of alternative haplotypes in heterozygous regions of the reference animal [24]. Marker placement on these segments suggests that these variants are tracking different haplotypes in the population, which is the desired intent of genetic markers used in genomic selection.

### Discussion

We have assembled a superior, extremely continuous reference assembly (Sscrofa11.1) by leveraging the excellent contig lengths provided by long reads, and a wealth of available data including Illumina paired-end, BAC end sequence, finished BAC sequence, fosmid end sequences, and the earlier curated draft assembly (Sscrofa10.2). The pig genome assemblies USMARCv1.0 and Sscrofa11.1 reported here are 92- and 694-fold, respectively, more continuous than the published draft reference genome sequence (Sscrofa10.2) [13]. The new pig reference genome assembly (Sscrofa11.1) with its contig N50 of 48,231,277 bp and 506 gaps compares favourably with the current human reference genome sequence (GRCh38.p13), which has a contig N50 of 57,879,411 bp and 875 gaps (Table 1). Indeed, considering only the chromosome assemblies built on PacBio long-read data (i.e., Sscrofa11—the autosomes SSC1-SSC18 plus SSCX), there are fewer gaps in the pig assembly than in human reference autosomes and HSAX assemblies. Most of the gaps in the Sscrofa11.1 reference assembly are attributed to the fragmented assembly of SSCY. The capturing of centromeres and telomeres for several chromosomes (Tables S2 and S3; Figs. S1 and S2) provides further evidence that the Sscrofa11.1 assembly is more complete. The increased contiguity of Sscrofa11.1 is evident in the graphical comparison with Sscrofa10.2 illustrated in Fig. 2.

The improvements in the reference genome sequence (Sscrofa11.1) relative to the draft assembly (Sscrofa10.2) [13] are not restricted to greater continuity and fewer gaps. The major flaws in the BAC clone-based draft assembly were (i) failures to resolve the sequence redundancy amongst sequence contigs within BAC clones and between adjacent overlapping BAC clones and (ii) failures to accurately order and orient the sequence contigs within BAC clones. Although the Sanger sequencing technology used has a much lower raw error rate than the PacBio





**Figure 2:** Visualization of improvements in assembly contiguity. Graphical visualization of contigs for Sscrofa11 (top) and Sscrofa10.2 (bottom) as alternating dark and light grey bars.

technology, the sequence coverage was only 4–6-fold across the genome. The improvements in continuity and quality (Table 2; Figs S5–S7) have yielded a better template for annotation, resulting in better gene models. The Sscrofa11.1 and USMARCv1.0 assemblies are classed as 4|4|1 and 3|5|1 ( $10^X$ : N50 contig [kb];  $10^Y$ : N50 scaffold [kb];  $Z = 1|0$ : assembled to chromosome level), respectively, compared to Sscrofa10.2 as 1|2|1 and the human GRCh38p5 assembly as 4|4|1 [45].

The improvement in the complete BUSCO genes indicates that both Sscrofa11.1 and USMARCv1.0 represent templates for annotation of gene models that are superior to the draft Sscrofa10.2 assembly and are comparable to the finished human and mouse reference genome sequences (Table S7). Furthermore, a companion bioinformatics analysis of available Iso-Seq and companion Illumina RNA-Seq data across the 9 tissues sur-

veyed has identified a large number (>54,000) of novel transcripts [36]. A majority of these transcripts are predicted to be spliced and validated by RNA-Seq data. Beiki and colleagues identified 10,465 genes expressing Iso-Seq transcripts that are present on the Sscrofa11.1 assembly but are unannotated in current NCBI or Ensembl annotations [36].

The alignment of the Sscrofa11.1 and USMARCv1.0 assemblies revealed that several of the USMARCv1.0 chromosome assemblies are inverted relative to Sscrofa11.1 and the cytogenetic map. Such inversions are due to the agnostic nature of genome assembly and post-assembly polishing programs. Unless these are corrected post hoc by manual curation, they result in artefactual inversions of the entire chromosome. However, such inversions do not generally affect downstream analysis that does not involve the relative order/orientation of whole chromosomes.

To ascertain whether the differences between Sscrofa11.1 and USMARCv1.0 in order and orientation within chromosomes represent assembly errors or real chromosomal differences will require further research. The sequence present at the telomeric end of the long arm of the USMARCv1.0 chromosome 7 assembly (after correcting the orientation of the USMARCv1.0 SSC7) is missing from the Sscrofa11.1 SSC7 assembly, and currently located on a 3.8-Mb unplaced scaffold (AEMK02000452.1). This unplaced scaffold harbours several genes including *DIO3*, *CKB*, and *NUDT14* whose orthologues map to human chromosome 14 as would be predicted from the pig-human comparative map [40]. This omission will be corrected in an updated assembly in future.

We demonstrate moderate improvements in the placement and ordering of commercial SNP genotyping markers on the Sscrofa11.1 reference genome that will affect future genomic selection programs. The reference-derived order of SNP markers plays a significant role in imputation accuracy, as demonstrated by a whole-genome survey of misassembled regions in cattle that found a correlation between imputation errors and misassemblies [46]. The gaps in SNP chip marker coverage that we identified will inform future marker selection surveys, which are likely to prioritize regions of the genome that are not currently being tracked by marker variants in close proximity to potential causal variant sites. In addition to the gaps in coverage provided by the commercial SNP chips there are regions of the genome assemblies that are devoid of annotated sequence variation as hitherto sequence variants have been discovered against incomplete genome assemblies. Thus, there is a need to re-analyse good-quality resequence data against the new assemblies in order to provide a better picture of sequence variation in the pig genome.

The cost of high-coverage whole-genome sequencing (WGS) precludes its routine use in breeding programs. However, it has been suggested that low coverage WGS followed by imputation of haplotypes may be a cost-effective replacement for SNP arrays in genomic selection [47]. Imputation from low coverage sequence data to whole-genome information has been shown to be highly accurate [48, 49]. At the 2018 World Congress on Genetics Applied to Livestock Production Aniek Bouwman reported that in a comparison of Sscrofa10.2 with Sscrofa11.1 (for SSC7 only) for imputation from 600 K SNP genotypes to whole-genome sequence, overall imputation accuracy on SSC7 improved considerably from 0.81 (1,019,754 variants) to 0.90 (1,129,045 variants) (A. Bouwman, personal communication). Thus, the improved assembly may not only serve as a better template for discovering genetic variation but also have advantages for genomic selection, including improved imputation accuracy.

Advances in the performance of long-read sequencing and scaffolding technologies, improvements in methods for assembling the sequence reads, and reductions in costs are enabling the acquisition of ever more complete genome sequences for multiple species and multiple individuals within a species. For example, in terms of adding species, the Vertebrate Genomes Project [50] aims to generate error-free, near-gapless, chromosomal-level, haplotyped phase assemblies of all of the ~66,000 vertebrate species and is currently in its first phase, which will see such assemblies created for an exemplar species from all 260 vertebrate orders. At the level of individuals within a species, smarter assembly algorithms and sequencing strategies are enabling the production of high quality truly haploid genome sequences for outbred individuals [24]. The establishment of assembled genome sequences for key individuals in the nucleus

populations of the leading pig breeding companies is achievable and potentially affordable. However, 10–30× genome coverage short-read data generated on the Illumina platform and aligned to a single reference genome is likely to remain the primary approach to sequencing multiple individuals within farmed animal species such as cattle and pigs [21, 51].

There are significant challenges in making multiple assembled genome resources useful and accessible. The current paradigm of presenting a reference genome as a linear representation of a haploid genome of a single individual is an inadequate reference for a species. As an interim solution the Ensembl team are annotating multiple assemblies for some species such as mouse and dog [52, 53, 54]. We have implemented this solution for pig genomes, including 11 Illumina short-read assemblies [17] in addition to the reference Sscrofa11.1 and USMARCv1.0 assemblies reported here (Ensembl release 98, September 2019) [39, 41]. Although these additional pig genomes are highly fragmented (Table S6) with contig N50 values from 32 to 102 kb, the genome annotation (Table S11) provides a resource to explore pig gene space across 13 genomes, including 6 Asian pig genomes. The latter are important given the deep phylogenetic split of ~1 million years between European and Asian pigs [13].

The current human genome reference already contains several hundred alternative haplotypes, and it is expected that the single linear reference genome of a species will be replaced with a new model—the graph genome [55–57]. These paradigm shifts in the representation of genomes present challenges for current sequence alignment tools and the “best-in-genome” annotations generated thus far. The generation of high quality annotation remains a labour-intensive and time-consuming enterprise. Comparisons with the human and mouse reference genome sequences, which have benefited from extensive manual annotation, indicate that there is further complexity in the porcine genome as yet unannotated (Table 3). It is very likely that there are many more transcripts, pseudogenes, and non-coding genes (especially long non-coding genes) to be discovered and annotated on the pig genome sequence [36]. The more highly continuous pig genome sequences reported here provide an improved framework against which to discover functional sequences, both coding and regulatory, and sequence variation. After correction for some contig/scaffold inversions in the USMARCv1.0 assembly, the overall agreement between the assemblies is high and illustrates that the majority of genomic variation is at smaller scales of structural variation. However, both assemblies still represent a composite of the 2 parental genomes present in the animals, with unknown effects of haplotype switching on the local accuracy across the assembly.

Future developments in top class genome sequences for the domestic pig are likely to include (i) gap closure of Sscrofa11.1 to yield an assembly with 1 contig per (autosomal) chromosome arm, exploiting the isogenic BAC and fosmid clone resource as illustrated here for chromosomes 16 and 18; and (ii) haplotype-resolved assemblies of a Meishan and White Composite F1 crossbred pig (i.e., the offspring of a Meishan sire and a White Composite dam that is approximately one-half Landrace, one-quarter Duroc, and one-quarter Yorkshire) currently being sequenced. Beyond this, haplotype-resolved assemblies for key genotypes in the leading pig breeding company nucleus populations and of miniature pig lines used in biomedical research can be anticipated in the next 5 years. Unfortunately, some of these genomes may not be released into the public domain. The first wave of results from the Functional Annotation of Animal

Genomes (FAANG) initiative [58, 59] are emerging and will add to the richness of pig genome annotation.

In conclusion, the new pig reference genome (Sscrofa11.1) described here represents a substantially enhanced resource for genetics and genomics research and applications for a species of importance to agriculture and biomedical research.

## Methods

Additional detailed methods and information on the assemblies and annotation are included in the Supplementary Materials.

### Preparation of genomic DNA

DNA was extracted from Duroc 2–14 cultured fibroblast cells passage 16–18 using the Qiagen Blood & Cell Culture DNA Maxi Kit. DNA was isolated from lung tissue from barrow MARC1423004 using a salt extraction method.

### Genome sequencing and assembly

Genomic DNAs from the samples described above were used to prepare libraries for sequencing on PacBio RS II sequencer (PacBio RS II Sequencing System, [RRID:SCR.017988](#)) [60]. For Duroc 2–14 DNA P6/C4 chemistry was used, whilst for MARC1423004 DNA a mix of P6/C4 and earlier P5/C3 chemistry was used.

Reads from the Duroc 2–14 DNA were assembled into contigs using the Falcon v0.4.0 assembly pipeline (Falcon, [RRID:SCR.016089](#)) following the standard protocol [27]. Quiver v. 2.3.0 [61] was used to correct the primary and alternative contigs. Only the primary pseudo-haplotype contigs were used in the assembly. The reads from the MARC1423004 DNA were assembled into contigs using Celera Assembler v8.3rc2 (Celera Assembler, [RRID:SCR.010750](#)) [32]. The contigs were scaffolded as described in the Results section.

### Fluorescence in situ hybridization

Metaphase preparations were fixed to slides and dehydrated through an ethanol series (2 mins each in 2× SSC, 70%, 85%, and 100% ethanol at room temperature). Probes were diluted in a formamide buffer (Cytocell) with Porcine Hybloc (Insight Biotech) and applied to the metaphase preparations on a 37°C hot plate before sealing with rubber cement. Probe and target DNA were simultaneously denatured for 2 mins on a 75°C hot plate prior to hybridization in a humidified chamber at 37°C for 16 h. Slides were washed after hybridization in 0.4× SSC at 72°C for 2 mins followed by 2× SSC/0.05% Tween 20 at room temperature for 30 sec, and then counterstained using VECTASHIELD anti-fade medium with DAPI (Vector Labs). Images were captured using an Olympus BX61 epifluorescence microscope with cooled CCD camera and SmartCapture (Digital Scientific UK) system.

### Analysis of repetitive sequences, including telomeres and centromeres

Repeats were identified using RepeatMasker (v.4.0.7, [RRID:SCR.012954](#)) [62] with a combined repeat database including Dfam (v.20170127) [63] and RepBase (v.20170127) [64]. RepeatMasker was run with “sensitive” (-s) setting using *sus scrofa* as the query species (-species “*sus scrofa*”). Repeats that showed >40% sequence divergence or were shorter than 70% of the expected sequence length were filtered out from subsequent analyses. The

presence of potentially novel repeats was assessed by RepeatMasker using the novel repeat library generated by RepeatModeler (v.1.0.11, [RRID:SCR.015027](#)) [62].

Telomeres were identified by running TRF [65] with default parameters apart from Mismatch (5) and Minscore (40). The identified repeat sequences were then searched for the occurrence of 5 identical, consecutive units of the TTAGGG vertebrate motif or its reverse complement and total occurrences of this motif were counted within the tandem repeat. Regions that contained ≥200 identical hexamer units, were >2 kb in length, and had a hexamer density of >0.5 were retained as potential telomeres.

Centromeres were predicted using the following strategy. First, the RepeatMasker output, both default and novel, was searched for centromeric repeat occurrences. Second, the assemblies were searched for known, experimentally verified, centromere-specific repeats [66, 67] in the Sscrofa11.1 genome. Then the 3 sets of repeat annotations were merged together with BEDTools (BEDTools, [RRID:SCR.006646](#)) [68] (median and mean length: 786 and 5,775 bp, respectively) and putative centromeric regions closer than 500 bp were collapsed into longer super-regions. Regions that were >5 kb were retained as potential centromeric sites.

### Long-read RNA sequencing (Iso-Seq)

The following tissues were harvested from MARC1423004 at age 48 days: brain (BioSamples: SAMN05952594), diaphragm (SAMN05952614), hypothalamus (SAMN05952595), liver (SAMN05952612), small intestine (SAMN05952615), skeletal muscle—longissimus dorsi (SAMN05952593), spleen (SAMN05952596), pituitary (SAMN05952626), and thymus (SAMN05952613). Total RNA from each of these tissues was extracted using Trizol reagent (ThermoFisher Scientific) and the provided protocol. Briefly, ~100 mg of tissue was ground in a mortar and pestle cooled with liquid nitrogen, and the powder was transferred to a tube with 1 mL of Trizol reagent added and mixed by vortexing. After 5 min at room temperature, 0.2 mL of chloroform was added and the mixture was shaken for 15 sec and left to stand another 3 min at room temperature. The tube was centrifuged at 12,000g for 15 min at 4°C. The RNA was precipitated from the aqueous phase with 0.5 mL of isopropanol. The RNA was further purified with extended DNaseI digestion to remove potential DNA contamination. The RNA quality was assessed with a Fragment Analyzer (Advanced Analytical Technologies Inc.). Only RNA samples of RQN > 7.0 were used for library construction. PacBio Iso-Seq libraries were constructed per the PacBio Iso-Seq protocol. Briefly, starting with 3 µg of total RNA, complementary DNA (cDNA) was synthesized by using the SMARTer PCR cDNA Synthesis Kit (Clontech) according to the Iso-Seq protocol (Pacific Biosciences). Then the cDNA was amplified using KAPA HiFi DNA Polymerase (KAPA Biotechnologies) for 10 or 12 cycles followed by purification and size selection into 4 fractions: 0.8–2, 2–3, 3–5, and >5 kb. The fragment size distribution was validated on a Fragment Analyzer (Advanced Analytical Technologies Inc.) and quantified on a DS-11 FX fluorometer (DeNovix). After a second round of large-scale PCR amplification and end repair, SMRT bell adapters were separately ligated to the cDNA fragments. Each size fraction was sequenced on 4 or 5 SMRT Cells v3 using P6-C4 chemistry and 6-h movies on a PacBio RS II sequencer (Pacific Biosciences). Short-read RNA-Seq libraries were also prepared for all 9 tissues using TruSeq stranded mRNA LT kits and supplied protocol (Illumina), and sequenced on an

Illumina NextSeq500 platform using v2 sequencing chemistry to generate  $2 \times 75$  bp paired-end reads.

The reads of interest were determined by using consensus-tools.sh in the SMRT-Analysis pipeline v2.0, with reads that were shorter than 300 bp and whose predicted accuracy was <75% removed. Full-length, non-concatemer (FLNC) reads were identified by running the classify.py command. The cDNA primer sequences as well as the poly(A) tails were trimmed prior to further analysis. Paired-end Illumina RNA-Seq reads from each tissue sample were trimmed to remove the adaptor sequences and low quality bases using Trimmomatic (v0.32, [RRID:SCR.011848](#)) [69] with explicit option settings: ILLUMINACLIP: adapters.fa: 2:30:10:1:true LEADING:3 TRAILING:3 SLIDINGWINDOW: 4:20 LEADING:3 TRAILING:3 MINLEN:25, and overlapping paired-end reads were merged using the PEAR software v0.9.6 (PEAR, [RRID:SCR.003776](#)) [70]. Subsequently, the merged and unmerged RNA-Seq reads from the same tissue samples were *in silico* normalized in a mode for single-end reads by using a Trinity v2.1.1 ([RRID:SCR.013048](#)) [71] utility, insilico.read.normalization.pl, with the following settings: -max\_cov 50 -max\_pct\_stdev 100 -single. Errors in the FLNC reads were corrected with the preprocessed RNA-Seq reads from the same tissue samples by using proovread (v2.12; Proovread, [RRID:SCR.017331](#)) [72]. Untrimmed sequences with at least some regions of high accuracy in the .trimmed.fq files were extracted based on sequence IDs in .untrimmed.fa files to balance off the contiguity and accuracy of the final reads.

### Short-read RNA sequencing

In addition to the Illumina short-read RNA-Seq data generated from MARC1423004 and used to correct the Iso-Seq data (see Long-read RNA sequencing (Iso-Seq) above), Illumina short-read RNA-Seq data (PRJEB19386) were also generated from a range of tissues from 4 juvenile Duroc pigs (2 male, 2 female) and used for annotation as described below. Extensive metadata with links to the protocols for sample collection and processing are linked to the BioSample entries under the Study Accession PRJEB19386. The tissues sampled are listed in Table S10. Sequencing libraries were prepared using a ribodepletion TruSeq stranded RNA protocol and 150-bp paired-end sequences generated on the Illumina HiSeq 2500 platform (Illumina HiSeq 2500 System, [RRID:SCR.016383](#)) in rapid mode.

### Annotation

The assembled genomes were annotated using the Ensembl pipelines (Ensembl, [RRID:SCR.002344](#)) [38] as detailed in the Supplementary Materials. The Iso-Seq and RNA-Seq data described above were used to build gene models.

### Mapping SNP chip probes

The probes from 4 commercial SNP chips were mapped to the Sscrofa10.2, Sscrofa11.1, and USMARCv1.0 assemblies using BWA MEM [73] and a wrapper script [74]. Probe sequence was derived from the marker manifest files that are available on the provider websites: Illumina PorcineSNP60 [1, 75], Affymetrix Axiom™ Porcine Genotyping Array [76], Gene Seek Genomic Profiler Porcine—HD beadChip [77], and Gene Seek Genomic Profiler Porcine v2—LD Chip [77]. To retain marker manifest coordinate information, each probe marker name was annotated with the chromosome and position of the marker's variant site from the manifest file. All mapping coordinates were tabulated into a single file and were sorted by the chromosome and position of

the manifest marker site. To derive and compare relative marker rank order, a custom Perl script [78] was used to sort and number markers based on their mapping locations in each assembly.

## Availability of Supporting Data and Materials

The genome assemblies are deposited at NCBI under accession numbers GCA.000003025.6 (Sscrofa11.1) and GCA.002844635.1 (USMARCv1.0). The associated BioSample accession numbers are SAMN02953785 and SAMN07325927, respectively. Iso-Seq and RNA-Seq data used for analysis and annotation are available under accession numbers PRJNA351265 and PRJEB19386, respectively. Supporting data and materials are available in the Giga-Science GigaDB database [79].

## Additional Files

### Supplementary Methods, Tables and Figures

**Table S1.** Pacific Biosciences read statistics.

**Table S2.** Predicted telomeres.

**Table S3.** Predicted centromeres.

**Table S4.** Assigning scaffolds to chromosomes.

**Table S5.** Alignment of radiation hybrid maps and genome assemblies.

**Table S6.** Assemblytics comparisons, assembly statistics.

**Table S7.** BUSCO results.

**Table S8.** Annotation statistics (Ensembl-NCBI comparison).

**Table S9.** Commercial SNP chip probes.

**Table S10.** Tissue samples.

**Table S11.** Ensembl annotation statistics for 13 pig genome assemblies.

**Figure S1.** Predicted telomeres.

**Figure S2.** Predicted centromeres.

**Figure S3.** Fluorescence *in situ* hybridization assignments.

**Figure S4.** Improvement in local order and orientation and reduction in redundancy.

**Figure S5.** Assembly comparisons in gEVAL (SSC15).

**Figure S6.** Assembly comparisons in gEVAL (SSC5).

**Figure S7.** Assembly comparisons in gEVAL (SSC18).

**Figure S8.** Order and orientation of SSC18 assemblies.

**Figure S9.** Order and orientation of SSC7 assemblies.

**Figure S10.** Order and orientation of SSC8 assemblies.

**Figure S11.** Assembly alignments.

**Figure S12.** Assemblytics results.

**Figure S13.** Counts of repetitive elements in 4 pig assemblies.

**Figure S14.** Average mapped length of repetitive elements in 4 pig genomes.

**Figure S15.** Assembly SNP rank concordance versus reported chromosomal location.

## Abbreviations

BAC: bacterial artificial chromosome; BLAST: Basic Local Alignment Search Tool; BLASTN: BLAST search of nucleotide database(s); bp: base pairs; BUSCO: Benchmarking Universal Single-Copy Orthologs; BWA: Burrows-Wheeler Aligner; CCD: charged couple device; cDNA: complementary DNA; DAPI: 4',6-diamidino-2-phenylindole; ENA: European Nucleotide Archive; FISH: fluorescence *in situ* hybridization; FLNC: full-length, non-concatemer; g: relative centrifugal force; Gb: gigabase pairs; GFF3: general feature format, version 3; GC: guanine-cytosine; HQ: high quality; ID: identity; Iso-Seq: long-read RNA sequencing using PacBio technology; kb: kilobase pairs; LC: low coverage;



LQ: low quality; LQLC: low quality, low coverage; MANE: Matched Annotation from NCBI and EMBL-EBI; Mb: megabase pairs; mRNA: messenger RNA; NCBI: National Center for Biotechnology Information; NIH: National Institutes of Health; PacBio: Pacific Biosciences; polyA: poly adenine; QC: quality control; RefSeq: NCBI Reference Sequence Database; RH: radiation hybrid; RNA-Seq: high-throughput short-read RNA sequencing; RQN: RNA quality number; RT: read through; SGSC: Swine Genome Sequencing Consortium; SMRT: single-molecule real-time; SNP: single-nucleotide polymorphism; SSC: saline sodium citrate; SSCn: *Sus scrofa* chromosome n; SD: standard deviation; TRF: Tandem Repeats Finder; USDA: United States Department of Agriculture; WGS: whole-genome sequencing.

## Competing Interests

R.H., K.K., and E.T. are employed by Pacific Biosciences; all other authors declare that they have no competing interests.

## Funding

This work was supported by the Biotechnology and Biological Sciences Research Council, Institute Strategic Programme Grant, BBS/E/D/20211550, A.L.A., M.W.; the Biotechnology and Biological Sciences Research Council, Institute Strategic Programme Grant, BBS/E/D/10002070, A.L.A., M.W.; the Biotechnology and Biological Sciences Research Council, Response Mode Grant, BB/F021372/1, N.A.; the Biotechnology and Biological Sciences Research Council, Response Mode Grant, BB/M011461/1, A.L.A.; the Biotechnology and Biological Sciences Research Council, Response Mode Grant, BB/M011615/1, P.F.; the Biotechnology and Biological Sciences Research Council, Response Mode Grant, BB/M01844X/1, A.L.A., M.W.; EU, FP7 Programme Quantomics, KBBE222664, A.L.A.; Wellcome Trust, WT108749/Z/15/Z, P.F.; USDA, CRIS Project, 8042-31000-001-00-D, D.M.B., B.D.R.; USDA, CRIS Project, 5090-31000-026-00-D, D.M.B.; USDA, CRIS Project, 3040-31000-100-00-D, T.P.L.S.

In addition to the funding acknowledged above we are grateful for support from the University of Cambridge, Department of Pathology, the European Molecular Biology Laboratory, and the Roslin Foundation. In addition H.L. and H.B. were supported by USDA NRSP-8 Swine Genome Coordination funding; S.K. and A.M.P. were supported by the Intramural Research Program of the National Human Genome Research Institute, US National Institutes of Health. This work used the computational resources of the NIH HPC Biowulf cluster (<https://hpc.nih.gov>) and the Iowa State University Lightning3 and ResearchIT clusters. The Ceres cluster (part of the USDA SCInet Initiative) was used to analyse part of this dataset.

## Authors' Contributions

A.L.A. and T.P.L.S. conceived, coordinated, and managed the project; A.L.A., P.F., D.A.H., T.P.L.S., and M.W. supervised staff and students performing the analyses; D.J.N., L.A.R., L.B.S., and T.P.L.S. provided biological resources; R.H., K.S.K., and T.P.L.S. generated PacBio sequence data; H.A.F., T.P.L.S., and R.T. generated Illumina WGS and RNA-Seq data; N.A.A., C.A.S., and B.M.S. provided SSCY assemblies; D.J.N. and T.P.L.S. generated Iso-Seq data; G.H., R.H., S.K., A.M.P., A.S.S., and A.W. generated sequence assemblies; A.W. polished and quality checked Sscrofa11.1; W.C., G.H., K.H., S.K., B.D.R., A.S.S., S.G.S., and E.T. performed quality checks on the sequence assemblies; R.E.O'C. and D.K.G. per-

formed cytogenetics analyses; L.E. analysed repeat sequences; H.B., H.L., N.M., and C.K.T. analysed Iso-Seq data; D.M.B. and G.A.R. analysed sequence variants; B.A., K.B., C.G.G., T.H., O.I., and F.J.M. annotated the assembled genome sequences; A.W. and A.L.A. drafted the manuscript; all authors read and approved the final manuscript.

## Acknowledgements

We are grateful to Chris Tyler-Smith (Wellcome Trust Sanger Institute) for sharing the SSCY sequence data for Sscrofa11.1.

## References

1. Ramos AM, Crooijmans RPMA, Affara NA, et al. Design of a high density SNP genotyping assay in the pig using SNPs identified and characterized by next generation sequencing technology. *PLoS One* 2009;4:e6524.
2. Hu ZL, Park CA, Reecy JM. Developmental progress and current status of the Animal QTLdb. *Nucleic Acids Res* 2016;44:D827–33.
3. Meuwissen T, Hayes B, Goddard M. Accelerating Improvement of livestock with genomic selection. *Annu Rev Anim Biosci* 2013;1:221–37.
4. Christensen OF, Madsen P, Nielsen B, et al. Single-step methods for genomic evaluation in pigs. *Animal* 2012;6:1565–71.
5. Cleveland M, Hickey JM. Practical implementation of cost-effective genomic selection in commercial pig breeding using imputation. *J Anim Sci* 2013;91:3583–92.
6. Vamathevan JJ, Hall MD, Hasan S, et al. Minipig and beagle animal model genomes aid species selection in pharmaceutical discovery and development. *Toxicol Appl Pharmacol* 2013;270:149–57.
7. Klymiuk N, Seeliger F, Bohlooly M, et al. Tailored pig models for preclinical efficacy and safety testing of targeted therapies. *Toxicol Pathol* 2016;44:346–57.
8. Wells KD, Prather RS. Genome-editing technologies to improve research, reproduction, and production in pigs. *Mol Reprod Dev* 2017;84:1012–7.
9. Servin B, Faraut T, Iannuccelli N, et al. High-resolution autosomal radiation hybrid maps of the pig genome and their contribution to the genome sequence assembly. *BMC Genomics* 2012;13:585.
10. Tortereau F, Servin B, Frantz L, et al. A high density recombination map of the pig reveals a correlation between sex-specific recombination and GC content. *BMC Genomics* 2012;13:586.
11. Yerle M, Lahbib-Mansais Y, Mellink C, et al. The PiGMap consortium cytogenetic map of the domestic pig (*Sus scrofa domestica*). *Mamm Genome* 1995;6:176–86.
12. Humphray SJ, Scott CE, Clark R, et al. A high utility integrated map of the pig genome. *Genome Biol* 2007;8:R139.
13. Groenen MAM, Archibald AL, Uenishi H, et al. Analyses of pig genomes provide insight into porcine demography and evolution. *Nature* 2012;491:393–8.
14. Warr A, Robert C, Hume D, et al. Identification of low-confidence regions in the pig reference genome (Sscrofa 10.2). *Front Genet* 2015;6:338.
15. O'Connor RE, Fonseka G, Frodsham R, et al. Isolation of subtelomeric sequences of porcine chromosomes for translocation screening reveals errors in the pig genome assembly. *Anim Genet* 2017;48:395–403.

16. Dawson HD, Chen C, Gaynor B, et al. The porcine translational research database: a manually curated, genomics and proteomics-based research resource. *BMC Genomics* 2017;18:643.
17. Li M, Chen L, Tian S, et al. Comprehensive variation discovery and recovery of missing sequence in the pig genome using multiple de novo assemblies. *Genome Res* 2017;27:865–74.
18. Schook LB, Beever JE, Rogers J, et al. Swine Genome Sequencing Consortium (SGSC): A strategic roadmap for sequencing the pig genome. *Comp Funct Genomics* 2005;6:251–5.
19. Robert C, Fuentes-Utrilla P, Troup K, et al. Design and development of exome capture sequencing for the domestic pig (*Sus scrofa*). *BMC Genomics* 2014;15:550.
20. Skinner BM, Sargent CA, Churcher C, et al. The pig X and Y chromosomes: Structure, sequence, and evolution. *Genome Res* 2016;26:130–9.
21. Frantz LAF, Schraiber JG, Madsen O, et al. Evidence of long-term gene flow and selection during domestication from analyses of Eurasian wild and domestic pig genomes. *Nat Genet* 2015;47:1141–8.
22. Groenen MAM. A decade of pig genome sequencing: A window on pig domestication and evolution. *Genet Sel Evol* 2016;48:23.
23. van Dijk EL, Jaszczyszyn Y, Naquin D, et al. The third revolution in sequencing technology. *Trends Genet* 2018;34:666–81.
24. Koren S, Rhie A, Walenz BP, et al. De novo assembly of haplotype-resolved genomes with trio binning. *Nat Biotechnol* 2018;36:1174–82.
25. CHORI-242: Porcine (*Sus scrofa*) BAC Library. <https://bacpacresources.org/library.php?id=124>. Accessed 17 April 2020.
26. Putnam NH, O'Connell BL, Stites JC, et al. Chromosome-scale shotgun assembly using an in vitro method for long-range linkage. *Genome Res* 2016;26:342–50.
27. Chin CS, Peluso P, Sedlazeck FJ, et al. Phased diploid genome assembly with single-molecule real-time sequencing. *Nat Methods* 2016;13:1050–4.
28. Kurtz S, Phillippy A, Delcher AL, et al. Versatile and open software for comparing large genomes. *Genome Biol* 2004;5:R12.
29. English AC, Richards S, Han Y, et al. Mind the gap: Upgrading genomes with Pacific Biosciences RS long-read sequencing technology. *PLoS One* 2012;7:e47768.
30. Altschul SF, Gish W, Miller W, et al. Basic Local Alignment Search Tool. *J Mol Biol* 1990;215:403–10.
31. Chow W, Brugger K, Caccamo M, et al. gEVAL—a web-based browser for evaluating genome assemblies. *Bioinformatics* 2016;32:2508–10.
32. Berlin K, Koren S, Chin CS, et al. Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. *Nat Biotechnol* 2015;33:623–30.
33. Nattestad M, Schatz MC. Assemblytics: A web analytics tool for the detection of variants from an assembly. *Bioinformatics* 2016;32:3021–3.
34. Simão FA, Waterhouse RM, Ioannidis P, et al. BUSCO: Assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 2015;31:3210–2.
35. Tseng E. CODING GENOME reconstruction Tool. 2017. <https://github.com/Magdoll/Cogent>. Accessed 7 September 2017.
36. Beiki H, Liu H, Manchanda N, et al. Improved annotation of the domestic pig genome through integration of Iso-Seq and RNA-seq data. *BMC Genomics* 2019;20:344.
37. Long Y, Su Y, Ai H, et al. A genome-wide association study of copy number variations with umbilical hernia in swine. *Anim Genet* 2016;47:298–305.
38. Cunningham F, Achuthan P, Akanni W, et al. Ensembl 2019. *Nucleic Acids Res* 2019;47(D1):D745–51.
39. Ensembl pig strains genome annotation Release 98, September 2019. <http://www.ensembl.org/Sus.scrofa/Info/Strains>. Accessed 15 October 2019.
40. Meyers SN, Rogatcheva MB, Larkin DM, et al. Piggy-BACing the human genome: II. A high-resolution, physically anchored, comparative map of the porcine autosomes. *Genomics* 2005;86:739–52.
41. Ensembl pig genome annotation Release 98, September 2019. <http://www.ensembl.org/Sus.scrofa/Info/Index>. Accessed 15 October 2019.
42. Ensembl archive. <http://may2017.archive.ensembl.org/Sus.scrofa/Info/Index>. Accessed 15 April 2020.
43. NCBI annotation report 106. [https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/Sus.scrofa/106/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Sus.scrofa/106/). Accessed 15 April 2020.
44. MANE (Matched Annotation from NCBI and EMBL-EBI). <https://www.ensembl.org/info/genome/genebuild/mane.html>. Accessed 15 April 2020.
45. gEVAL: Genome Evaluation Browser. <https://geval.sanger.ac.uk/>. Accessed 15 October 2019.
46. Utsunomiya ATH, Santos DJ, Boison SA, et al. Revealing mis-assembled segments in the bovine reference genome by high resolution linkage disequilibrium scan. *BMC Genomics* 2016;17:705.
47. Hickey JM. Sequencing millions of animals for genomic selection 2.0. *J Anim Breed Genet* 2013;130:331–2.
48. Le SQ, Durbin R. SNP detection and genotyping from low-coverage sequencing data on multiple diploid samples. *Genome Res* 2011;21:952–60.
49. Li Y, Sidore C, Kang HM, et al. Low-coverage sequencing: Implications for design of complex trait association studies. *Genome Res* 2011;21:940–51.
50. Vertebrate Genomes Project. <https://vertebrategenomesproject.org/>. Accessed 15 October 2019.
51. Daetwyler HD, Capitan A, Pausch H, et al. Whole-genome sequencing of 234 bulls facilitates mapping of monogenic and complex traits in cattle. *Nat Genet* 2014;46:858–65.
52. Lilue J, Doran AG, Fiddes IT, et al. Sixteen diverse laboratory mouse reference genomes define strain-specific haplotypes and novel functional loci. *Nat Genet* 2018;50:1574–83.
53. Ensembl mouse strain annotation. <http://www.ensembl.org/Mus.musculus/Info/Strains>. Accessed 15 April 2020.
54. Ensembl dog breed annotation. <http://www.ensembl.org/Canis.familiaris/Info/Strains>. Accessed 15 April 2020.
55. Baier U, Beller T, Ohlebusch E. Graphical pan-genome analysis with compressed suffix trees and the Burrows-Wheeler transform. *Bioinformatics* 2015;32:497–504.
56. Chaisson MJP, Wilson RK, Eichler EE. Genetic variation and the de novo assembly of human genomes. *Nat Rev Genet* 2015;16:627–40.
57. Garrison E, Sirén J, Novak AM, et al. Variation graph toolkit improves read mapping by representing genetic variation in the reference. *Nat Biotechnol* 2018;36:875–9.
58. Andersson L, Archibald AL, Bottema CD, et al. Coordinated international action to accelerate genome-to-phenome with FAANG, the Functional Annotation of Animal Genomes project. *Genome Biol* 2015;16:57.
59. Foissac S, Djebali S, Munyard K, et al. Multispecies annotation of transcriptome and chromatin structure in domesticated animals. *BMC Biol* 2019;17:108.

60. Pendleton M, Sebra R, Pang AA, et al. Assembly and diploid architecture of an individual human genome via single-molecule technologies. *Nat Methods* 2015;12: 780–6.
61. Chin CS, Alexander DH, Marks P, et al. Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nat Methods* 2013;10:563–9.
62. RepeatMasker. <http://www.repeatmasker.org>. Accessed 17 April 2020.
63. Hubley R, Finn RD, Clements J, et al. The Dfam database of repetitive DNA families. *Nucleic Acids Res* 2016;44: D81–9.
64. Bao W, Kojima KK, Kohany O. Repbase Update, a database of repetitive elements in eukaryotic genomes. *Mob DNA* 2015;6(1), doi:10.1186/s13100-015-0041-9.
65. Benson G. Tandem Repeats Finder: A program to analyze DNA sequences. *Nucleic Acids Res* 1999;27: 573–80.
66. Miller JR, Hindkjær J, Thomsen PD. A chromosomal basis for the differential organization of a porcine centromere-specific repeat. *Cytogenet Cell Genet* 1993;62: 37–41.
67. Riquet J, Mulsant P, Yerle M, et al. Sequence analysis and genetic mapping of porcine chromosome 11 centromeric S0048 marker. *Cytogenet Cell Genet* 1996;74:127–32.
68. Quinlan AR, Hall IM. BEDTools: A flexible suite of utilities for comparing genomic features, *Bioinformatics* 2010;26: 841–2.
69. Bolger AM, Lohse M, Usadel B. Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics* 2014;30:2114–20.
70. Zhang J, Kobert K, Flouri T, et al. PEAR: A fast and accurate Illumina Paired-End reAd mergeR. *Bioinformatics* 2014;30:614–20.
71. Grabherr MG, Haas BJ, Yassour M, et al. Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat Biotechnol* 2011;29:644–52.
72. Hackl T, Hedrich R, Schultz J, et al. Proovread: Large-scale high-accuracy PacBio correction through iterative short read consensus. *Bioinformatics* 2014;30:3004–11.
73. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 2009;25:1754–60.
74. Bickhart D. Align and order SNP probes. 2019. [https://github.com/njdbickhart/perl.toolchain/blob/master/assembly\\_scripts/alignAndOrderSnPProbes.pl](https://github.com/njdbickhart/perl.toolchain/blob/master/assembly_scripts/alignAndOrderSnPProbes.pl).
75. Illumina PorcineSNP60 Beadchip. <https://emea.illumina.com/products/by-type/microarray-kits/porcine-snp60.html>. Accessed 17 April 2020.
76. Axiom™ Porcine Genotyping Array. <https://www.thermofisher.com/order/catalog/product/550588#/550588>. Accessed 17 April 2020.
77. Gene Seek Genomic Profiler Porcine. <https://genomics.neogen.com/uk/ggp-porcine>. Accessed 17 April 2020.
78. Bickhart D. Pig genome SNP sort rank order. 2019. [https://github.com/njdbickhart/perl.toolchain/blob/master/assembly\\_scripts/pigGenomeSNPSortRankOrder.pl](https://github.com/njdbickhart/perl.toolchain/blob/master/assembly_scripts/pigGenomeSNPSortRankOrder.pl). Accessed 11 May 2019.
79. Warr A, Affara N, Aken B, et al. Supporting data for “An improved pig reference genome sequence to enable pig genetics and genomics research.” GigaScience Database 2020. <http://dx.doi.org/10.5524/100732>.